

Date: _____

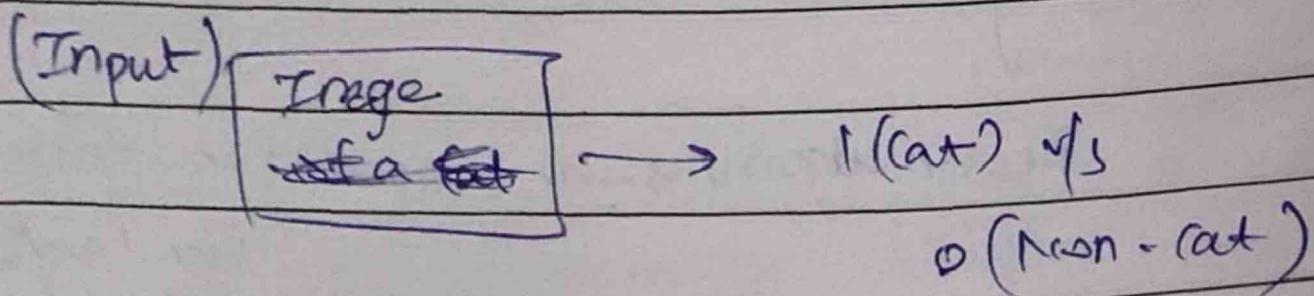
WEEK 2

Neural Networks Basics

→ Binary classification.

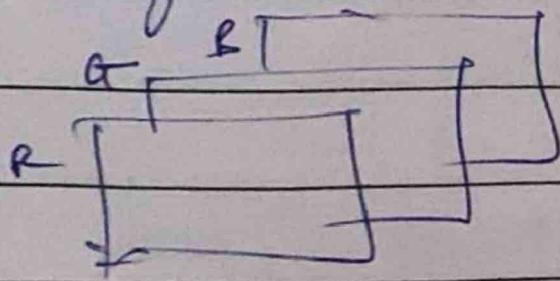
logistic regression is an algorithm for
binary classification.

Example:



$$x \rightarrow y$$

Image:



(matrix of pixels)

3 matrices
of 64×64

dimension

Date:

We unroll all these pixel intensity values
int \in a feature vector x

$$x = \begin{matrix} 126 \\ 25 \\ \vdots \\ 255 \\ 312 \\ 306 \end{matrix} \quad \begin{matrix} 64 \times 64 \times 3 \\ = 12288 \end{matrix}$$

Dimension of

input feature
vector = $n = n_x$

$$= 12288$$

$$x \rightarrow y$$

Notation :-

$$(x, y) \quad x \in \mathbb{R}^{n_x}, y \in \{0, 1\}$$

m training examples : $(x^{(1)}, y^{(1)}) (x^{(2)}, y^{(2)}) \dots (x^{(m)}, y^{(m)})$

M_{train}
training set

M_{test}

test set

→ Input matrix

Date: _____

$$X = \begin{bmatrix} ; & ; & ; \\ ; & \{ & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ ; & ; & \} & | \\ & & & n_x \end{bmatrix} \uparrow \downarrow$$

$\leftarrow m \rightarrow$

$$\underline{X \in \mathbb{R}^{n_x \times m}}$$

$$X.\text{shape} = (n_x, m)$$

→ output matrix

$$Y = [y^{(1)} \quad y^{(2)} \quad \dots \quad y^{(m)}]$$

$$\underline{Y \in \mathbb{R}^{1 \times m}}$$

$$Y.\text{shape} = (1, m)$$

Logistic Regression

Date:

Given x , want $\hat{y} = P(y=1|X)$

$$x \in \mathbb{R}^{n_x}$$

Parameters: $w \in \mathbb{R}^{n_x}$, $b \in \mathbb{R}$

Output $\hat{y} = w^T x + b$

For this algo formula, for y being 1
 \hat{y} has to be $0 \leq \hat{y} \leq 1$.

$\therefore w^T x$ can be $>> 1$

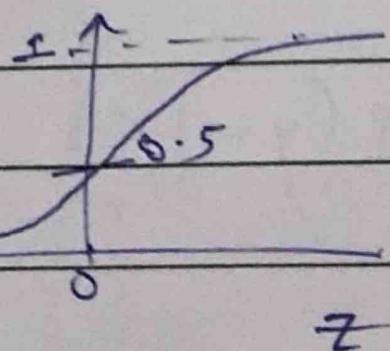
(or) it can even be negative.

$$\therefore \hat{y} = \sigma(w^T x + b)$$

sigmoid function

Date:

$\sigma(z)$



$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

If z is very large

$$\therefore \sigma(z) \approx \frac{1}{0+1} = 1$$

If z is large negative number

$$\sigma(z) = \frac{1}{1 + e^{-z}} \approx 0$$

Date:

Logistic Regression Cost function :

$$\hat{y} = \sigma(w^T x + b), \text{ where } \sigma(z) = \frac{1}{1 + e^{-z}}$$

Given $\{(x^{(1)}, y^{(1)}); \dots; (x^{(m)}, y^{(m)})\}$

want $\hat{y}^{(i)} \approx y^{(i)}$

Loss (Error) function :

$$L(\hat{y}, y) = \frac{1}{2} (\hat{y} - y)^2$$

(Not used)

This formula causes
optimisation problem.

→ loss-function is defined for a single training example

Date:

$$L(\hat{y}, y) = - (y \log \hat{y} + (1-y) \log (1-\hat{y}))$$

If $y=1$: $L(\hat{y}, y) = -\log \hat{y}$

∴ $-\log \hat{y}$ should be ↓

∴ $\log \hat{y}$ ↑

∴ \hat{y} ↑.

If $y=0$: $L(\hat{y}, y) = -\log (1-\hat{y})$

∴ $-\log (1-\hat{y})$ ↓

∴ $\log (1-\hat{y})$ ↑

∴ \hat{y} ↓

→ Cost function measures over an entire training set.

Date:

Cost function: $\mathcal{J}(w, b)$

$$= \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

$$= -\frac{1}{m} \sum_{i=1}^m \left[y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log (1-\hat{y}^{(i)}) \right]$$

Gradient Descent:

Recap:

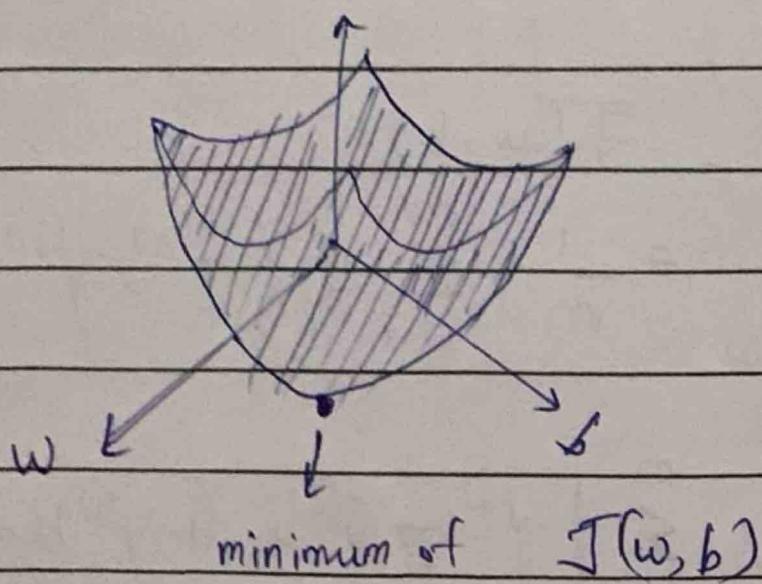
$$\hat{y} = \sigma(w^T x + b), \quad \sigma(z) = \frac{1}{1+e^{-z}}$$

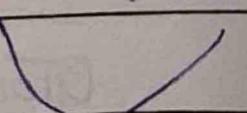
$$\mathcal{J}(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

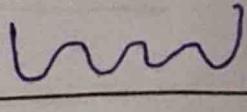
$$= -\frac{1}{m} \sum_{i=1}^m \left[y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log (1-\hat{y}^{(i)}) \right]$$

Want to find w, b so that
 $\mathcal{J}(w, b)$ is minimum.

Date:



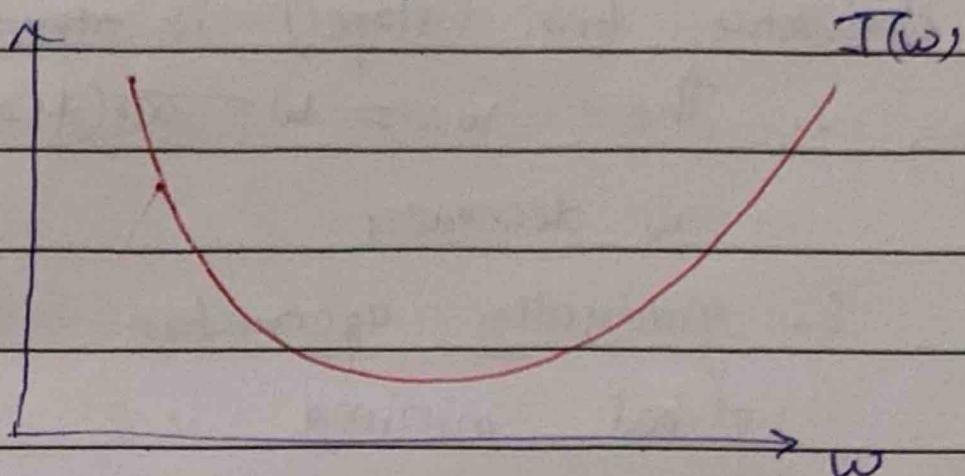
$J(w, b)$ [Convex function] 

[Non-convex functions] 

To find good values of parameters,
we initialise values of w and b ,
most often with values of zero.

Date:

Gradient descent graph:



Repeat {

$$w := w - \alpha \frac{d J(w)}{dw}$$

}

→ α : learning rate (controls how big a step we take on each iteration of gradient descent)

→ $\frac{d J(w)}{dw}$ will be written as "dw"

is the update of the change you want to make to the parameters w .

Date:

on RHS of the graph;

the derivative term (slope) is ~~+~~
positive, \therefore the $w := w - \alpha \cdot (\text{the})$
 w decreases

- gradually approaches
global minima.

On LHS of the graph;

the derivative term (slope) is

negative, $\therefore w := w - \alpha \cdot (-\text{ve})$
 $\therefore w$ increases

- gradually approaches
global minima.

\rightarrow for $J(w, b)$

$$w := w - \alpha \frac{dJ(w, b)}{dw}$$

$$\frac{\delta J(w, b)}{\delta w}$$

$$b := b - \alpha \frac{dJ(w, b)}{db}$$

$$\frac{\delta J(w, b)}{\delta b}$$

Date: 22/07/2024

Computation Graph:

example:

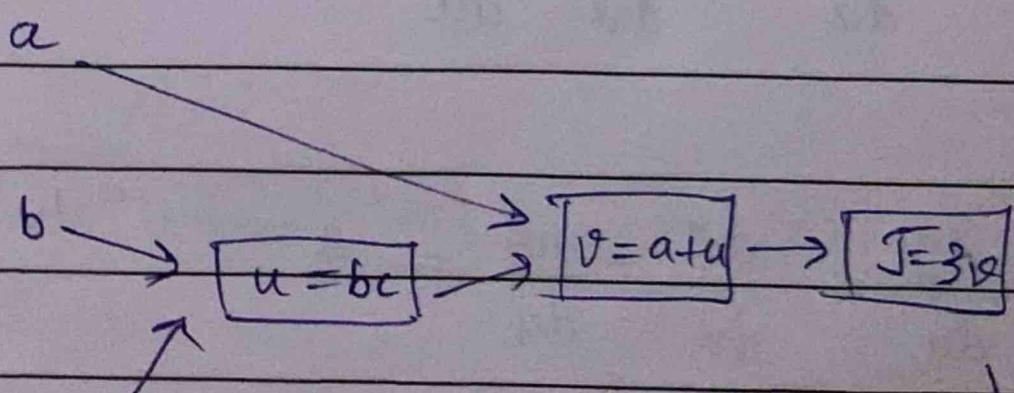
$$J(a, b, c) = 3(a + bc)$$

→ 3 steps:

① $u = bc$

② $v = a + u$

③ $J = 3v$



Computation graph :

This is the

Final output

variable which will

be used for optimisation.

Date: _____

For example:

If we want to compute derivative of the
~~last step~~ output variable, then we're
done one step of backward propagation.
~~called~~ one step backward in graph.

$$\frac{dJ}{dv} = 3 \quad \text{"d}\nu\text{"}$$

$$\frac{dJ}{da} = \frac{dJ}{dv} \cdot \frac{dv}{da} = 3 \times 1 = 3 \quad \text{"d}\alpha\text{"}$$

$$\frac{dJ}{du} = \frac{dJ}{dv} \cdot \frac{dv}{du} = 3 \quad \text{"d}\upsilon\text{"}$$

$$\frac{dJ}{db} = \frac{dJ}{dv} \cdot \frac{dv}{du} \cdot \frac{du}{db} = 3c = 3(2) \\ = 6$$

"d}\beta\text{"}

Date:

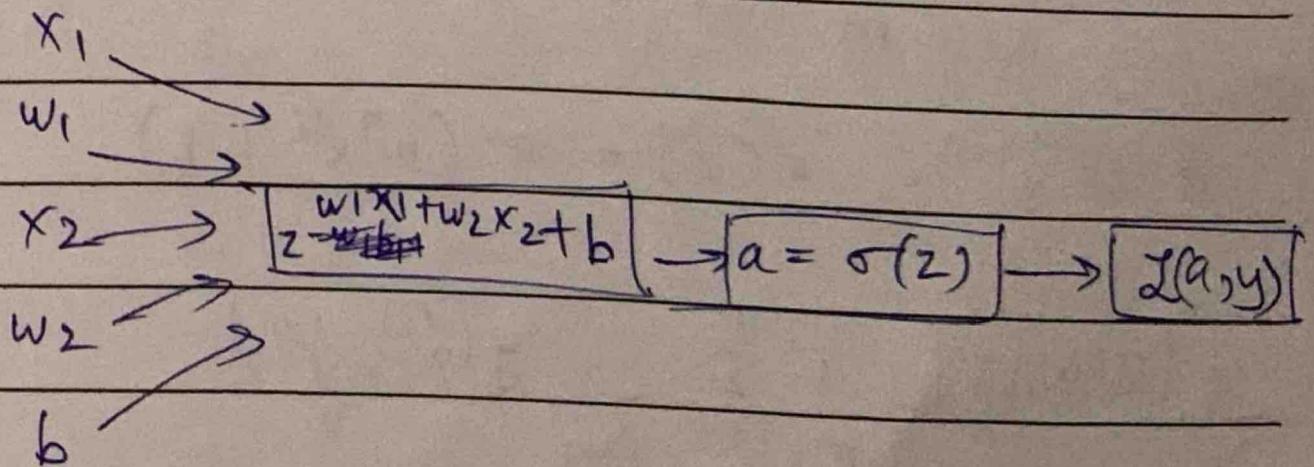
To compute derivatives of all the variables, follow right-to-left computation (backward propagation).

logistic regression recap:

$$z = w^T x + b$$

$$\hat{y} = a = \sigma(z)$$

$$\mathcal{L}(a, y) = -(y \log(a) + (1-y) \log(1-a))$$



$$\frac{da}{dz} = \frac{-y}{a} + \frac{(1-y)}{1-a}$$

$$\frac{d\mathcal{L}}{dz} = \frac{da}{dz} \cdot \frac{d\mathcal{L}}{da} = \frac{da}{dz} - a(1-a)$$

$$\therefore \frac{d\mathcal{L}}{dz} = a - y$$

Date:

$$\text{"d}w_1\text{"} = x_1 \cdot dz$$

$$\text{"d}w_2\text{"} = x_2 \cdot dz$$

$$\text{"d}b\text{"} = dz$$

Logistic Regression on 'm' examples =

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(a^{(i)}, y^{(i)})$$

$$a^{(i)} = \hat{y}^{(i)} = \sigma(z^{(i)}) = \sigma(w^T x^{(i)} + b)$$

$$\frac{\partial [J(w, b)]}{\partial w_1} = \frac{1}{m} \sum_{i=1}^m \frac{\partial \mathcal{L}(a^{(i)}, y^{(i)})}{\partial w_1}$$

{①, ② = for loops }

Date:

$$J=0, dw_1=0, dw_2=0, db=0$$

①

for $i=1$ to m {

$$z^{(i)} = w^T x^{(i)} + b \quad dw = np.zeros(n_x, 1)$$

$$a^{(i)} = \sigma(z^{(i)})$$

$$J+ = -[y^{(i)} \log a^{(i)} + (1-y^{(i)}) \log(1-a^{(i)})]$$

$$d_z^{(i)} = a^{(i)} - y^{(i)}$$

$$\left. \begin{array}{l} dw_1+ = x^{(i)} d_z^{(i)} \\ dw_2+ = x_2^{(i)} d_z^{(i)} \end{array} \right\} \begin{array}{l} \text{assuming } n_x=2 \\ \text{2 features.} \end{array}$$

$$db+ = d_z^{(i)}$$

vectorised:

$$dw+ = X^{(i)} d_z^{(i)}$$

$$|J|=m$$

$$|dw_1|=m; |dw_2|=m; |db|=m$$

}

Note: Here we are

$w_1 := w_1 - \alpha dw_1$ using two for loops

$w_2 := w_2 - \alpha dw_2$ making it inefficient,

$b := b - \alpha db$ hence vectorisation can be used.

Date:

Vectorisation:

→ To get rid of explicit for loops in our code, to make it more efficient.

$$z = w^T x + b$$

$$w \in \mathbb{R}^{hn}$$

$$w = \begin{bmatrix} ; \\ ; \\ ; \end{bmatrix} \quad x = \begin{bmatrix} ; \\ ; \\ ; \end{bmatrix} \quad x \in \mathbb{R}^{nx}$$

Non-vectorised:

$$z = 0$$

for i in range(n-n) : {

$$z += w[i] * x[i] \}$$

$$z += b$$

Date:

Vectorised:

$$z = \underbrace{\text{np.dot}(w, x)}_{w^T \cdot x} + b$$

GPUs } use SIMD - single instruction
CPUs } multiple data
to process in-built functions of
python.

vectors and matrix valued functions :

$$v = \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{bmatrix} \rightarrow u = \begin{bmatrix} e^{v_1} \\ e^{v_2} \\ \vdots \\ e^{v_n} \end{bmatrix}$$

Non-vectorised:

$u = \text{np.zeros}((n, 1))$

for i in range(n):

$u[i] = \text{math.exp}(v[i]);$

vectorised:

`import numpy as np`

$u = \text{np.exp}(v)$

$\text{np.log}(v)$

$\text{np.abs}(v)$

$\text{np.maximum}(v, 0)$

Date:

Vectorizing Logistic Regression

$$z^{(i)} = w^T \cdot x^{(i)} + b \quad z^{(m)} = w^T \cdot x^{(m)} + b$$

.....

$$a^{(i)} = \sigma(z^{(i)}) \quad a^{(m)} = \sigma(z^{(m)})$$

$$X = \begin{bmatrix} 1 & 1 & \dots & 1 \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ 1 & 1 & \dots & 1 \end{bmatrix} \quad \begin{matrix} (n_n, m) \\ n_n \times m \\ \mathbb{R} \end{matrix}$$

$Z = \downarrow$

$$\begin{bmatrix} z^{(1)} \\ z^{(2)} \\ \vdots \\ z^{(m)} \end{bmatrix} = w^T X + \begin{bmatrix} b \\ b \\ b \\ b \end{bmatrix} \quad \begin{matrix} (1 \times n_n) \\ (n_n \times m) \\ = 1 \times m \end{matrix}$$

$$= \begin{bmatrix} w^T x^{(1)} + b \\ z^{(1)} \\ w^T x^{(2)} + b \\ z^{(2)} \\ \vdots \\ \vdots \\ w^T x^{(m)} + b \\ z^{(m)} \end{bmatrix} \quad \begin{matrix} 1 \times m \end{matrix}$$

$$Z = np.\text{dot}(w^T, x) + b$$

$$A = [a^{(1)} \ a^{(2)} \ \dots \ a^{(m)}] = \sigma(Z)$$

Date:

$$dz^{(1)} = a^{(1)} - y^{(1)} \quad dz^{(2)} = a^{(2)} - y^{(2)}$$

— —

$$dZ = \begin{bmatrix} dz^{(1)} & dz^{(2)} & \dots & dz^{(m)} \end{bmatrix}$$

$1 \times m$

$$A = \begin{bmatrix} a^{(1)} & a^{(2)} & \dots & a^{(m)} \end{bmatrix}$$

$$Y = \begin{bmatrix} y^{(1)} & y^{(2)} & \dots & y^{(m)} \end{bmatrix}$$

$$dZ = A - Y = \begin{bmatrix} a^{(1)} - y^{(1)} & a^{(2)} - y^{(2)} & \dots & a^{(m)} - y^{(m)} \end{bmatrix}$$

$$dw = 0$$

$$dw+ = x^{(1)} dz^{(1)}$$

$$dw+ = x^{(2)} dz^{(2)}$$

}

$$db = 0$$

$$db+ = dz^{(1)}$$

$$db+ = dz^{(2)}$$

{ } db+ = dz^{(m)}

$$dw/m = m$$

$$db/m = m$$

$$db = \frac{1}{m} \sum_{i=1}^m dz^{(i)} = \frac{1}{m} [np \cdot \text{sum}(dZ)]$$

$$dw = \frac{1}{m} (X \cdot dZ^T)$$

Date: _____

$$Z = w^T \cdot X + b$$

$$= np.\text{dot}(w.T, X) + b$$

$$A = \sigma(Z)$$

$$dZ = A - Y$$

$$dw = \frac{1}{m} [(X \cdot dZ^T)]$$

$$db = \frac{1}{m} [np.\text{sum}(dZ)]$$

$$w := w - \alpha dw$$

$$b := b - \alpha db$$

Date:

Broadcasting Example :

	Apples	Beets	Eggs	Potatoes
carb				
Protein				
Fat				

$$cal = A \cdot \text{sum}(axis=0)$$

$$\text{percentage} = 100 * A / \text{cal. reshape}(1, 4)$$

General Principle :

$$\begin{matrix} \cdot & (m, n) & + & (l, n) \\ \text{matrix} & \nearrow & \swarrow & \end{matrix} \longrightarrow (m, n)$$
$$(m, l) \longrightarrow (m, n)$$

$$\begin{bmatrix} m_{1,1} \\ 2 \\ 3 \\ 1 \end{bmatrix} + R = \begin{bmatrix} 101 \\ 102 \\ 103 \\ 1 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 2 & 3 & \dots \end{bmatrix} + 100 = \begin{bmatrix} 101 & 102 & 103 & \dots \end{bmatrix}$$